

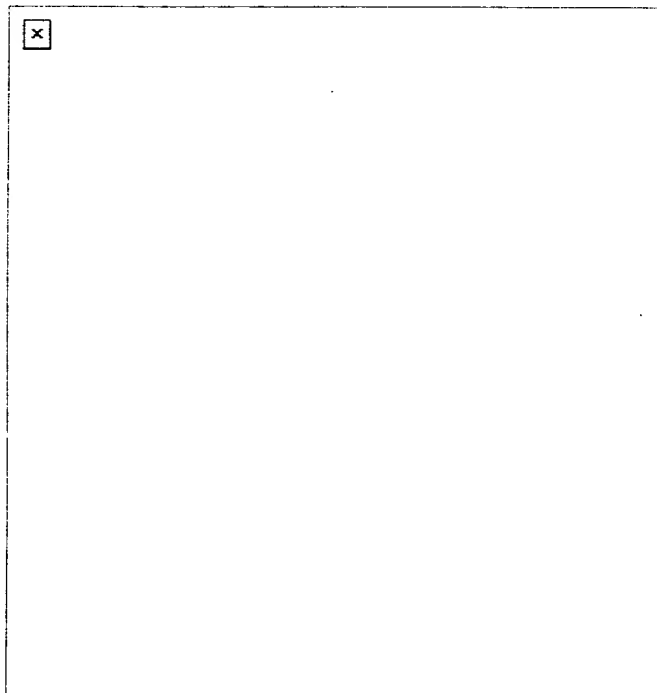
SYSTEM AND METHOD FOR FULL-TEXT RETRIEVAL AND RECORD MEDIUM WHERE FULL-TEXT RETRIEVING PROGRAM IS RECORDED

Patent number: JP10312395
Publication date: 1998-11-24
Inventor: KITAGAWA RYOKO; SHIRAI TADAHIRO; SUZUKI YOSHIAKI; SUGIYAMA SHINYA; SUGAYA TOMOHIDE
Applicant: TOSHIBA CORP
Classification:
- international: G06F17/30
- european:
Application number: JP19970324462 19971126
Priority number(s):

Abstract of JP10312395

PROBLEM TO BE SOLVED: To improve the retrieval precision by reducing retrieval noise and to reduce the necessary storage capacity of a data base.

SOLUTION: The full-text retrieval system has plural kinds of retrieval information tables 18 to 22 that specify documents which are set to the different numbers of characters on character units and include specified character units, generates plural kinds of character unit groups from an inputted retrieval key word, selects the retrieval information tables by the generated character unit groups, and performs retrieval from the selected retrieval information tables in the generated character units to specify a document including the retrieval key word among respective documents including the retrieved character units. In this case, the number of retrieval information tables is set to a different value by the kinds of characters of a character string including character units set in respective retrieval information tables, and the kind of characters of the retrieval key word is decided to select a retrieval information table to be retrieved according to the character kind and the number of characters of the character string.



Data supplied from the *esp@cenet* database - Worldwide

(19)日本国特許庁 (J P)

(12) 公 開 特 許 公 報 (A)

(11)特許出願公開番号

特開平10-312395

(43)公開日 平成10年(1998)11月24日

(51)Int.Cl.⁶

G 0 6 F 17/30

識別記号

F I

G 0 6 F 15/40

15/411

15/413

3 7 0 A

3 1 0

3 1 0 A

3 1 0 B

審査請求 未請求 請求項の数7 O L (全 15 頁)

(21)出願番号 特願平9-324462

(22)出願日 平成9年(1997)11月26日

(31)優先権主張番号 特願平9-54729

(32)優先日 平9(1997)3月10日

(33)優先権主張国 日本(J P)

(71)出願人 000003078

株式会社東芝

神奈川県川崎市幸区堀川町72番地

(72)発明者 北川 良子

東京都府中市東芝町1番地 株式会社東芝
府中工場内

(72)発明者 白井 直裕

東京都府中市東芝町1番地 株式会社東芝
府中工場内

(72)発明者 鈴木 善昭

東京都府中市東芝町1番地 株式会社東芝
府中工場内

(74)代理人 弁理士 鈴江 武彦 (外6名)

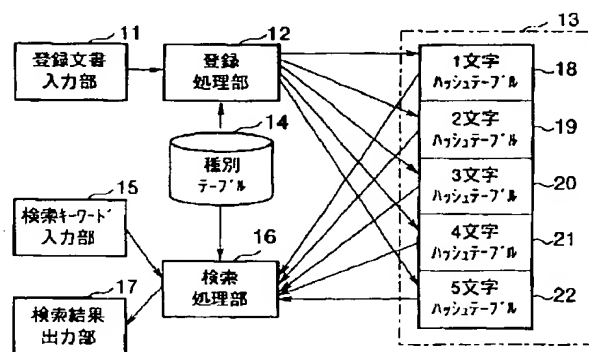
最終頁に続く

(54)【発明の名称】 全文検索システム及び全文検索方法並びに全文検索プログラムを記録した記録媒体

(57)【要約】

【課題】 検索ノイズを低減して検索精度を向上できると共に、データベースの必要記憶容量を低減する。

【解決手段】 それぞれ文字単位の文字数が異なる値に設定され、指定された文字単位が含まれる文書を特定する複数種類の検索情報テーブル18~22を有し、入力された検索キーワードから複数種類の文字単位群を生成し、この生成された各文字単位群毎に検索情報テーブルを選択して、この選択した各検索情報テーブルを生成された各文字単位で検索し、検索された各文字単位が含まれる各文書から検索キーワードが含まれる文書を特定する全文検索システムにおいて、検索情報テーブルの設定数を、各検索情報テーブルに設定される文字単位を含む文字列の文字種別毎に異なる値に設定し、かつ検索キーワードの文字種別を判定して、この文字種別と文字列の文字数とから検索すべき検索情報テーブルを選択する。



【特許請求の範囲】

【請求項1】 それぞれ文字単位の文字数が異なる値に設定され、指定された文字単位が含まれる文書を特定する複数種類の検索情報テーブルを有し、入力された検索キーワードから複数種類の文字単位群を生成し、この生成された各文字単位群毎に検索情報テーブルを選択して、この選択した各検索情報テーブルを前記生成された各文字単位で検索し、検索された各文字単位が含まれる各文書から前記検索キーワードが含まれる文書を特定する全文検索システムにおいて、前記検索情報テーブルの設定数は、各検索情報テーブルに設定される文字単位を含む文字列の文字種別毎に異なる値を有し、

前記検索キーワードの文字種別を判定して、この文字種別と文字列の文字数とから前記検索すべき検索情報テーブルを選択することを特徴とする全文検索システム。

【請求項2】 それぞれ文字単位の文字数が異なる値に設定され、指定された文字単位が含まれる文書を特定する複数種類の検索情報テーブルを有し、入力された検索キーワードから複数種類の文字単位群を生成し、この生成された各文字単位群毎に検索情報テーブルを選択して、この選択した各検索情報テーブルを前記生成された各文字単位で検索し、検索された各文字単位が含まれる各文書から前記検索キーワードが含まれる文書を特定する全文検索システムにおいて、

前記各文字単位は、この文字単位を構成する1個又は複数の文字からハッシュ関数を用いて算出されたハッシュ値で示され、

前記複数種類の検索情報テーブルには、文字数が互いに異なる複数種類の文字単位に対応するハッシュ値が組込まれた共通検索情報テーブルが含まれることを特徴とする全文検索システム。

【請求項3】 それぞれ文字単位の文字数が異なる値に設定され、指定された文字単位が含まれる文書を特定する複数種類の検索情報テーブルを有し、入力された検索キーワードから複数種類の文字単位群を生成し、この生成された各文字単位群毎に検索情報テーブルを選択して、この選択した各検索情報テーブルを前記生成された各文字単位で検索し、検索された各文字単位が含まれる各文書から前記検索キーワードが含まれる文書を特定する全文検索システムにおいて、

前記各文字単位は、この文字単位を構成する1個又は複数の文字からハッシュ関数を用いて算出されたハッシュ値で示され、

前記複数種類の検索情報テーブルには、前記文字単位を含む文字列の文字種別毎に異なるハッシュ関数を用いて算出されたハッシュ値が設定された検索情報テーブルが含まれ、

前記検索キーワードの文字種別を判定して、文字数と文字種別とから前記検索すべき検索情報テーブルを選択す

ることを特徴とする全文検索システム。

【請求項4】 それぞれ文字単位の文字数が異なる値に設定され、指定された文字単位が含まれる文書を特定する複数種類の検索情報テーブルを有し、入力された検索キーワードから複数種類の文字単位群を生成し、この生成された各文字単位群毎に検索情報テーブルを選択して、この選択した各検索情報テーブルを前記生成された各文字単位で検索し、検索された各文字単位が含まれる各文書から前記検索キーワードが含まれる文書を特定する全文検索方法において、

前記検索情報テーブルの設定数を、各検索情報テーブルに設定される文字単位を含む文字列の文字種別毎に異なる値に設定し、

前記検索キーワードの文字種別を判定して、この文字種別と文字列の文字数とから前記検索すべき検索情報テーブルを選択する全文検索方法。

【請求項5】 それぞれ文字単位の文字数が異なる値に設定され、指定された文字単位が含まれる文書を特定する複数種類の検索情報テーブルを有し、入力された検索キーワードから複数種類の文字単位群を生成し、この生成された各文字単位群毎に検索情報テーブルを選択して、この選択した各検索情報テーブルを前記生成された各文字単位で検索し、検索された各文字単位が含まれる各文書から前記検索キーワードが含まれる文書を特定する全文検索システムにおける全文検索プログラムを記録したコンピュータ読取り可能な記録媒体であって、

前記検索キーワードの設定数を、各検索情報テーブルに設定される文字単位を含む文字列の文字種別毎に異なる値とさせ、

前記検索キーワードの文字種別を判定させ、この文字種別と文字列の文字数とから前記検索すべき検索情報テーブルを選択させることを特徴とする全文検索プログラムを記録したコンピュータ読取り可能な記録媒体。

【請求項6】 それぞれ文字単位の文字数が異なる値に設定され、指定された文字単位が含まれる文書を特定する複数種類の検索情報テーブルを有し、入力された検索キーワードから複数種類の文字単位群を生成し、この生成された各文字単位群毎に検索情報テーブルを選択して、この選択した各検索情報テーブルを前記生成された各文字単位で検索し、検索された各文字単位が含まれる各文書から前記検索キーワードが含まれる文書を特定する全文検索システムにおける全文検索プログラムを記録したコンピュータ読取り可能な記録媒体であって、

前記各文字単位を、この文字単位を構成する1個又は複数の文字からハッシュ関数を用いて算出されたハッシュ値で示させ、

前記複数種類の検索情報テーブルには、文字数が互いに異なる複数種類の文字単位に対応するハッシュ値が組み込まれた共通検索情報テーブルを含ませることを特徴とする全文検索プログラムを記録したコンピュータ読取り

可能な記録媒体。

【請求項7】 それぞれ文字単位の文字数が異なる値に設定され、指定された文字単位が含まれる文書を特定する複数種類の検索情報テーブルを有し、入力された検索キーワードから複数種類の文字単位群を生成し、この生成された各文字単位群毎に検索情報テーブルを選択して、この選択した各検索情報テーブルを前記生成された各文字単位で検索し、検索された各文字単位が含まれる各文書から前記検索キーワードが含まれる文書を特定する全文検索システムにおける全文検索プログラムを記録したコンピュータ読取り可能な記録媒体であって、前記各文字単位をこの文字単位を構成する1個又は複数の文字からハッシュ関数を用いて算出されたハッシュ値で示させ、前記複数種類の検索情報テーブルには、前記文字単位を含む文字列の文字種別毎に異なるハッシュ関数を用いて算出されたハッシュ値が設定された検索情報テーブルを含ませ、前記検索キーワードの文字種別を判定させ、文字数と文字種別とから前記検索すべき検索情報テーブルを選択させることを特徴とする全文検索プログラムを記録したコンピュータ読取り可能な記録媒体。

【発明の詳細な説明】

【0001】

【発明の属する技術分野】 本発明は、大量の文書から必要とする情報が記載された文書を検索する情報検索システムに係わり、特に比較的長い検索キーワードで必要な文書を検索できる全文検索システム、全文検索方法、及び全文検索プログラムを記録した記録媒体に関する。

【0002】

【従来の技術】 一般の情報検索システムのデータベースに新規の文書（文献）を登録する場合は、該当文書に含まれる複数のキーワードをデータベースに登録する。このキーワードは一般に予め決められた単語で構成されている。そして、この情報検索システムを用いて必要な情報が記載された文書（文献）を調べる場合は、必要な情報に関係するキーワードでデータベースを検索すると、このキーワードが登録された文書の文書名、発行所、著者、発行日、所蔵場所等の該当文書を特定する情報が検索結果として出力される。

【0003】 このような情報検索システムにおいては、付与したキーワードでしか検索できなかった。この不都合を解消するために文書中の任意の単語や文字列でデータベースを検索できる全文検索システムが開発されている。

【0004】 図12は全文検索システムの概略構成図である。この全文検索システムは、大きく分けて、検索キーワードが入力されるキーボード等の検索キーワード入力部1と、実際に検索を実行する検索処理部2と、データベース3と、検索結果を出力するCRT表示装置や印字装置等で構成された検索結果出力部4とで構成されて

いる。

【0005】 データベース3内には、例えば、1文字検索情報テーブル5、2文字検索情報テーブル6、3文字検索情報テーブル7等の複数の検索情報テーブルが設けられている。

【0006】 1文字検索情報テーブル5内には、図13(a)に示すように、ひらかな、カタカナ、漢字、英字、数字等の各1文字単位毎に、該当1文字単位がこのデータベース3に登録された各文書に含まれるか否かの情報が記憶されている。具体的には、図示するように、該当1文字単位が含まれる文書番号#に[1]のビットが設定され、該当1文字が含まれない文書番号#に[0]のビットが設定される。

【0007】 2文字検索情報テーブル6内には、図13(b)に示すように、前記ひらかな、カタカナ、漢字、英字、数字の2文字の全ての組合わせを示す2文字単位毎に、該当2文字単位が含まれる文書番号#に[1]のビットが設定され、該当2文字単位が含まれない文書番号#に[0]のビットが設定される。したがって、2文字検索情報テーブル6に設定されている2文字単位の数は1文字検索情報テーブル5に設定されている1文字からなる1文字単位の数のほぼ2乗値となる。

【0008】 3文字検索情報テーブル7内には、図示しないが、ひらかな、カタカナ、漢字、英字、数字の3文字の全ての組合わせを示す3文字単位毎に、該当3文字単位が含まれる文書番号#に[1]のビットが設定され、該当1文字が含まれない文書番号#に[0]のビットが設定される。

【0009】 そして、検索処理部2は図14に示す流れ図に従って検索キーワード入力部1から入力された検索キーワードに対する検索処理を実行する。流れ図のステップS1において、検索キーワード入力部1から一つの検索キーワードが入力されると、S2において、該当検索キーワードを1文字単位に分解する。例えば、図15に示すように、検索キーワードが[日本の技術]のように2つの単語と1つの助詞で構成された5文字の文字列からなる場合は、[日][本][の][技][術]のように1文字からなる5つの1文字単位に分割する。S3にて、5つの各1文字単位に対してそれぞれ1文字検索情報テーブル5を検索して、各1文字単位が含まれる各文書番号#を得る。

【0010】 S4にて、入力された検索キーワードが2文字以上で構成されていた場合は、該当検索キーワードを2文字単位に分解する。例えば、検索キーワードが[日本の技術]の場合は、図15に示すように、[日本][本の][の技][技術]の4つの2文字単位に分割する(S5)。そして、4つの各2文字単位に対してそれぞれ2文字検索情報テーブル6を検索して、各2文字単位が含まれる各文書番号#を得る(S6)。

【0011】 S7にて、入力された検索キーワードが3

文字以上の文字列で構成されていた場合は、該当検索キーワードを3文字単位に分解する。例えば、検索キーワードが「日本の技術」の場合は、図15に示すように、「日本の」「本の技」「の技術」の3つの3文字単位に分割する(S8)。そして、3つの各3文字単位に対してそれぞれ3文字検索情報テーブル7を検索して、各3文字単位が含まれる各文書番号#を得る(S9)。

【0012】そして、S10において、上述したS3、S6、S9にて実行された全ての検索結果のAND(AND)結果を得る。具体的には、全ての文字単位が含まれる文書番号#を抽出する。

【0013】例えば、検索キーワードが「日本の技術」の場合は、1文字単位と2文字単位と3文字単位との合計12個の文字単位が全て含まれる文書番号#を抽出して、この文書番号#をこの検索処理部2における検索結果として検索結果出力部4を介して出力する(S11)。

【0014】このようにして、全文検索システムにおいては、2つ以上の単語を含む比較的長い文字列からなる検索キーワードで該当文字列を含む文書をデータベース3から検索できる。

【0015】このように構成された全文検索システムにおいては、長い文字列が含まれる文書を精度よく検索するためには、1、2、3文字検索情報テーブル5、6、7以外にも4文字、5文字の検索情報テーブルが必要である。

【0016】しかし、検索情報テーブルに設定する文字単位の文字数が増加すると、組合わせ数が飛躍的に増加し、データベース3が必要とする記憶容量が大幅に増大する。

【0017】このような不都合を解消するために、データベース3の各検索情報テーブル5、6、7に登録されている1個又は複数の文字からなる文字単位をハッシュ関数 $F_i(c_1, c_2, \dots, c_i)$ を用いて算出されたハッシュ値 h で示す技術が開発されている。

【0018】図16は、文字単位の構成文字数が i である場合における各文字単位をハッシュ値 $h(=1, 2, 3, \dots, N_i)$ で示した場合の検索情報テーブルである i 文字ハッシュテーブル8を示す。

【0019】このハッシュ関数 $F_i(c_1, c_2, \dots, c_i)$ は、各文字単位を構成する i 個の各文字 c_1, c_2, \dots, c_i の関数で示される。したがって、この関数 F_i の式を調整することによって、各文字単位の各文字 c_1, c_2, \dots, c_i の複数種類の組合わせに対して同一ハッシュ値 h を設定可能である。

【0020】その結果、 i 文字数で構成される文字単位の全ての組合わせの数より、ハッシュ値 h の数を大幅に減少できる。よって、このハッシュ関数 $F_i(c_1, c_2, \dots, c_i)$ を用いることによって、各検索情報テーブルの記憶容量をある程度低減できる。

【0021】

【発明が解決しようとする課題】しかしながら、上述したハッシュ関数 $F_i(c_1, c_2, \dots, c_i)$ を用いた全文検索システムにおいても、まだ解消すべき次のような課題があった。すなわち、前述したように、文字数の大きい文字列を精度よく検索するには、高い文字数の検索情報テーブル(文字ハッシュテーブル)を設ける必要があるが、データベース3の記憶容量の制約からせいぜい図12に示したように、1、2、3文字検索情報テーブル(文字ハッシュテーブル)までである。

【0022】したがって、図15に示したような5文字からなる文字列「日本の技術」を検索する場合は、5文字の文字列そのもので検索情報テーブルを検索していないので、実際には目的の文字列が含まれない文書が検索されてしまうこともあり、検索精度が低下する。なお、検索結果に混入するこの誤った文書を「検索ノイズ」と称する。

【0023】検索精度が低下するのは、文字種の多い日本語の漢字よりも文字種の少ない英字、数字でよく発生することが知られている。文字種が少ない場合は、検索する文字列の並びが文書中に多数出現する場合が多い。例えば、0～9の数値では、3文字の組合せが $10 \times 10 \times 10 = 10^3$ 通りしかない。

【0024】しかしJIS規格の第一水準のかな漢字は約3000個存在するので、3文字の組合せが $3000 \times 3000 \times 3000 = 2.7 \times 10^{10}$ 通りあるので、それぞれの出現頻度が低くなる。

【0025】例えば、5文字からなる「10000」の数値を検索する場合、各検索情報テーブルに対して次の検索が行われる。

1文字テーブル 1, 0, 0, 0, 0
2文字テーブル 10, 00, 00, 00
3文字テーブル 100, 000, 000

次に、4文字からなる「1000」の数値を検索することを考えると、各検索情報テーブルに対して次の検索が行われる。

【0026】1文字テーブル 1, 0, 0, 0
2文字テーブル 10, 00, 00
3文字テーブル 100, 000

この5文字の文字列「10000」の検索と4文字の文字列「1000」の検索とは、同じ文字単位の検索を複数回行うため、実際には両者ともに次に示す全く同じ検索を行うことになる。

【0027】1文字テーブル 1, 0
2文字テーブル 10, 00
3文字テーブル 100, 000

すなわち、出現頻度の高い「00」の文字単位を何度も検索することになり、たとえ各テーブルにおけるそれぞれの検索結果のANDをとったとしても、同一の検索結果が出力される懸念があり、検索精度が低下する。

【0028】そのため、文字列「10000」を含む文書を検索しても、文字列「1000」を含む文書が同時に検索されてしまう。また、5文字の文字列「12000」を検索する場合は、次の文字単位を検索することになる。

【0029】1文字テーブル 1, 2, 0, 0, 0

2文字テーブル 12, 20, 00, 00

3文字テーブル 120, 200, 000

そのため、次のような文書も誤って検索してしまう。

【0030】「2000件の中で120件が・・・」

このように、検索する文字列の文字種が少ない場合は、検索する文字列の文字並びが文書中に多数出現する場合が多く、前述した検索ノイズが増加する傾向にある。

【0031】そこで4文字、5文字を検索するための4文字検索情報テーブル、5文字検索情報テーブルを増設することが考えられるが、データベース3の必要記憶容量が飛躍的に増大し、実用的でないという問題が発生する

【0032】また、データベース3の必要記憶容量を低減させるために、前述した図16に示すハッシュ関数 $F_i(c_1, c_2, \dots, c_i)$ を採用したハッシュテーブル8を採用することが考えられる。

【0033】しかし、ハッシュ関数 $F_i(c_1, c_2, \dots, c_i)$ を使用する場合は、前述したように、必然的に、異なる文字単位が同じハッシュ値 h を持つ可能性があるため、同じハッシュ値 h を持つ別の文字単位を含む文書が検索ノイズとして検索される懸念がある。

【0034】特にハッシュテーブル8中のばらつきに片寄りがある場合は、検索ノイズが増加する。一般的に、このハッシュ関数 $F_i(c_1, c_2, \dots, c_i)$ の設定は難しく、ハッシュ関数 F_i の設定の精度によって検索ノイズの発生率が増減する。

【0035】このように、従来の全文検索システムにおいては、検索する際には、検索キーワードを1文字単位、2文字単位、3文字単位ずつ区切った組合わせて検索するため、4文字以上の文字列の場合、特に数値などの文字種別において、正しい検索結果が得られない懸念がある。

【0036】また、高い検索精度を得るために文字単位に含まれる文字数が多い多数の検索情報テーブルを組込むことが考えられるが、データベース3の必要記憶容量が大幅に増加する問題があった。さらに、ハッシュ関数 F_i を用いてデータベース3の必要記憶容量を減少させる場合は、ハッシュ関数 F_i の設定の精度によって検索ノイズの発生率が増加する懸念がある。

【0037】本発明はこのような事情に鑑みてなされたものであり、検索情報テーブルの設定数を検索すべき文字列の文字種別に応じた値に設定することによって、検索情報テーブルを記憶するデータベースの記憶容量を大幅に増加することなく、検索精度を向上できる全文検索

システム及び全文検索方法並びに全文検索プログラムを記録した記録媒体を提供することを目的とする。

【0038】また、各検索情報テーブルの構成を工夫することによって、検索精度を低下することなく、検索情報テーブルを記憶するデータベースの記憶容量を低減できる全文検索システムを提供することを目的とする。

【0039】

【課題を解決するための手段】本発明は、それぞれ文字単位の文字数が異なる値に設定され、指定された文字単位が含まれる文書を特定する複数種類の検索情報テーブルを有し、入力された検索キーワードから複数種類の文字単位群を生成し、この生成された各文字単位群毎に検索情報テーブルを選択して、この選択した各検索情報テーブルを生成された各文字単位で検索し、検索された各文字単位が含まれる各文書から検索キーワードが含まれる文書を特定する全文検索システムに適用される。

【0040】そして、上記課題を解消するために請求項1においては、検索情報テーブルの設定数は、各検索情報テーブルに設定される文字単位を含む文字列の文字種別毎に異なる値を有し、検索キーワードの文字種別を判定して、この文字種別と文字列の文字数とから前記検索すべき検索情報テーブルを選択する。

【0041】また、請求項2においては、各文字単位は、この文字単位を構成する1個又は複数の文字からハッシュ関数を用いて算出されたハッシュ値で示され、複数種類の検索情報テーブルには、文字数が互いに異なる複数種類の文字単位に対応するハッシュ値が組込まれた共通検索情報テーブルが含まれる。

【0042】請求項3においては、各文字単位は、この文字単位を構成する1個又は複数の文字からハッシュ関数を用いて算出されたハッシュ値で示され、複数種類の検索情報テーブルには、文字単位を含む文字列の文字種別毎に異なるハッシュ関数を用いて算出されたハッシュ値が設定された検索情報テーブルが含まれ、検索キーワードの文字種別を判定して、文字数と文字種別とから記検索すべき検索情報テーブルを選択する。

【0043】請求項4においては、それぞれ文字単位の文字数が異なる値に設定され、指定された文字単位が含まれる文書を特定する複数種類の検索情報テーブルを有し、入力された検索キーワードから複数種類の文字単位群を生成し、この生成された各文字単位群毎に検索情報テーブルを選択して、この選択した各検索情報テーブルを生成された各文字単位で検索し、検索された各文字単位が含まれる各文書から検索キーワードが含まれる文書を特定する全文検索方法において、検索情報テーブルの設定数を、各検索情報テーブルに設定される文字単位を含む文字列の文字種別毎に異なる値に設定し、検索キーワードの文字種別を判定して、この文字種別と文字列の文字数とから検索すべき検索情報テーブルを選択する。

【0044】このように構成された全文検索システム及

び全文検索方法においては、検索情報テーブルの設定数は、各検索情報テーブルに設定される文字単位を含む文字列の文字種別毎に異なる値を有している。

【0045】すなわち、文字種別として、例えば英字、数字、ひらがな、カタカナ、漢字第一水準、漢字第二水準、外字等が存在する。そして、前述したように、検索すべき文字列が英字、数字のみで構成されていた場合は、たとえ構成文字数が異なる文字列であったとしても同一文字列を検索する事態が発生する確率が多いので、検索ノイズが発生する確率が高くなる。

【0046】一方、英字、数字以外のひらがな、カタカナ、漢字第一水準、漢字第二水準、外字等の文字種別においては、文字列に同一文字が多数含まれることは希であるので、少ない文字数の文字単位の検索情報テーブルのみを用いて検索したとしても検索ノイズの発生確率は少ない。

【0047】したがって、検索ノイズの発生確率が高い英数字で構成される文字列を分離した文字単位に対する検索情報テーブルに対してのみ、構成文字数の高い検索情報テーブルを設ければ、検索ノイズの発生確率が低下する。

【0048】なお、英字、数字の文字種は漢字に比較して格段に少ないので、例えば4文字や5文字の検索情報テーブルを設けたとしても、各検索情報テーブルを記憶するデータベースの記憶容量が大幅に増加することはない。

【0049】また、別の発明においては、各検索情報テーブルに登録される文字単位は、この文字単位を構成する1個又は複数の文字からハッシュ関数を用いて算出されたハッシュ値で示されている。

【0050】このハッシュ値は各文字数毎の検索情報テーブル毎に設定されるが、たとえ構成文字数が異なる文字単位であっても、発生確率の小さい文字単位どうしであれば、このハッシュ値を一つの共通検索情報テーブルに組込んだとしても検索ノイズの発生確率は大幅に上昇することはない。

【0051】よって、検索精度を低下させることなく、データベースの必要記憶容量を節減できる。また、別の発明においては、複数種類の検索情報テーブルには、文字単位を含む文字列の文字種別毎に異なるハッシュ関数を用いて算出されたハッシュ値が設定されている。

【0052】すなわち、検索対象となる文字列の文字種別は、前述したように、英字、数字、ひらがな、カタカナ、漢字第一水準、漢字第二水準、外字等が存在する。一般に、複数の文字からなる文字単位を例えば文書から無作為に抽出した場合は、検出される文字単位の各文字の各組み合わせの発生確率は文字種別に応じてそれぞれ異なる。

【0053】したがって、発生確率の高い文字単位どうしを同一ハッシュ値に設定されることを防止すると、各

文字種別毎にハッシュ関数を設定することによって、検索ノイズの発生確率を低減できる。

【0054】また、英字、数字は文字種が少ないので、同一組合せに対してできるだけ複数のハッシュ値が設定されないようにハッシュ関数を設定すればよい。さらに、請求項5乃至請求項7の発明は、それぞれ文字単位の文字数が異なる値に設定され、指定された文字単位が含まれる文書を特定する複数種類の検索情報テーブルを有し、入力された検索キーワードから複数種類の文字単位群を生成し、この生成された各文字単位群毎に検索情報テーブルを選択して、この選択した各検索情報テーブルを生成された各文字単位で検索し、検索された各文字単位が含まれる各文書から前記検索キーワードが含まれる文書を特定する全文検索システムにおける全文検索プログラムを記録したコンピュータ読取り可能な記録媒体である。

【0055】そして、請求項5における全文検索プログラムは、コンピュータに対して、検索キーワードの設定数を、各検索情報テーブルに設定される文字単位を含む文字列の文字種別毎に異なる値とさせ、検索キーワードの文字種別を判定させ、この文字種別と文字列の文字数とから検索すべき検索情報テーブルを選択させる。

【0056】また、請求項6における全文検索プログラムは、コンピュータに対して、各文字単位を、この文字単位を構成する1個又は複数の文字からハッシュ関数を用いて算出されたハッシュ値で示させ、複数種類の件策情報テーブルには、文字数が互いに異なる複数種類の文字単位に対応するハッシュ値が組み込まれた共通検索情報テーブルを含ませる。

【0057】さらに、請求項7における全文検索プログラムは、コンピュータに対して、各文字単位をこの文字単位を構成する1個又は複数の文字からハッシュ関数を用いて算出されたハッシュ値で示させ、複数種類の検索情報テーブルには、文字単位を含む文字列の文字種別毎に異なるハッシュ関数を用いて算出されたハッシュ値が設定された検索情報テーブルを含ませ、検索キーワードの文字種別を判定させ、文字数と文字種別とから前記検索すべき検索情報テーブルを選択させる。

【0058】このように構成された全文検索プログラムを記録した記録媒体を用いることによって、上述した機能を有していない従来の全文検索システムに対しても、簡単に上述した機能を付加することができる。

【0059】

【発明の実施の形態】以下本発明の各実施形態を図面を用いて説明する。

(第1実施形態)図1は本発明の第1実施形態の全文検索方法を適用した全文検索システムの概略構成を示すブロック図である。

【0060】この全文検索システムは、大きく分けて、文章をデータベース13へ登録するため登録文書入力部

11及び登録処理部12と、データベース13から必要な情報を検索するための検索キーワード入力部15、検索処理部16及び検索結果出力部17と、種別テーブル14とで構成されている。

【0061】データベース13内には、図示するように、複数種類の検索情報テーブルとしての1文字ハッシュテーブル18、2文字ハッシュテーブル19、3文字ハッシュテーブル20、4文字ハッシュテーブル21及び5文字ハッシュテーブル22の合計5つハッシュテーブルが設けられている。

【0062】各ハッシュテーブル18～22内には、各ハッシュテーブル18～22で指定されて1文字～5文字の各文字単位を構成する文字のすべての組み合わせをハッシュ関数 F_1 (c_1, c_2, \dots, c_i)を用いて算出されたハッシュ値 h が設定されており、該当ハッシュ値 h が得られる文字単位を含む文書番号 $\#$ が[1]のビットで登録されている。

【0063】この合計5つハッシュテーブル18～22に登録される1文字～5文字の各組み合わせの文字単位を含む文字列の文字種別として、この第1実施形態システムにおいては、

- (a) 英字、数字のみ
- (b) ひらがな、カタカナ、漢字第一水準を含む
- (c) 漢字第二水準、外字のみ

の3種類に区分している。

【0064】1文字ハッシュテーブル18内には、図3(a)に示すように、(a)、(b)、(c)の3種類のすべての文字種別を含む4000文字種の各1文字(1文字単位)毎に、該当1文字単位からハッシュ関数 F_1 (c_1)を用いて算出されたハッシュ値 h ($=1 \sim N_1$)が設定されている。

【0065】2文字ハッシュテーブル19内には、図3(b)に示すように、(a)、(b)、(c)の3種類のすべての文字種別を含む約4000文字種の2文字単位毎に、該当2文字単位からハッシュ関数 F_2 (c_1, c_2)を用いて算出されたハッシュ値 h ($=1 \sim N_2$)が設定されている。

【0066】3文字ハッシュテーブル20内には、図3(c)に示すように、(c)の漢字第二水準、外字を除く、(a)、(b)の2種類の文字種別を含む約3500文字種の3文字単位毎に、該当3文字単位からハッシュ関数 F_3 (c_1, c_2, c_3)を用いて算出されたハッシュ値 h ($=1 \sim N_3$)が設定されている。

【0067】4文字ハッシュテーブル21内には、図4(a)に示すように、(a)の英字、数字の1種類のみの文字種別を含む約50文字種の4文字単位毎に、該当4文字単位からハッシュ関数 F_4 (c_1, c_2, c_3, c_4)を用いて算出されたハッシュ値 h ($=1 \sim N_4$)が設定されている。

【0068】5文字ハッシュテーブル22内には、図4

(b)に示すように、(a)の英字、数字の1種類のみの文字種別を含む約50文字種の5文字単位毎に、該当5文字単位からハッシュ関数 F_5 (c_1, c_2, c_3, c_4, c_5)を用いて算出されたハッシュ値 h ($=1 \sim N_5$)が設定されている。

【0069】前記種別テーブル14内には、図2に示すように、検索キーワード、登録文字列を構成する前述した(a)(b)(c)の3種類の文字種別毎に、検索対象又は登録対象の各ハッシュテーブル18～22が登録されている。

【0070】具体的には、(a)の英字、数字に対しては全てのハッシュテーブル18～22が登録され、(b)のひらがな、カタカナ、漢字第一水準に対しては1、2、3文字のハッシュテーブル18～20が登録され、(c)の漢字第二水準、外字に対しては1、2文字のハッシュテーブル18、19のみが登録されている。

【0071】次に、登録処理部12が行うデータベース13の各ハッシュテーブル18～22に対する登録文書入力部11から入力された文書の登録処理を図5に示す流れ図を用いて説明する。

【0072】ステップR1において、データベース13に対して登録すべき文書が存在することを確認すると、該当文書を読取る(R2)。そして、この文書内に登録すべき文字列が存在すると(R3)、該当文字列を読み込み(R4)、この入力した文字列の文字数 K 、及び該当文字列の文字種別を判断する。具体的には、前述した(a)、(b)、(c)に区分する(R5)。

【0073】そして、種別テーブル14から判別された文字種別に対応する使用ハッシュテーブルを特定する(R6)。特定された使用ハッシュテーブルの数 n と、使用ハッシュテーブル名を $Na(1)$ 、 $Na(2)$ 、 $Na(3)$ 、 \dots 、 $Na(n)$ と設定する(R7)。

【0074】以上の準備処理が終了すると、使用ハッシュテーブルを特定するインデックス i を1に初期化する(R8)。そして、インデックス i が使用ハッシュテーブル数 n 以下で、かつ使用ハッシュテーブル名 $Na(i)$ が文字列の文字数 K 以下の場合(R9)、 i 文字ハッシュテーブルに対する文書番号 $\#$ の登録処理を開始する。

【0075】すなわち、該当文字列を i 個の連続文字からなる複数の文字単位に分割して(R10)、この各文字単位からハッシュ関数 F_i (c_1, c_2, \dots, c_i)を用いて各ハッシュ値 h を算出して、各ハッシュテーブル18～22のうちの i 文字ハッシュテーブルの該当ハッシュ値 h の欄に対して該当文書番号 $\#$ を登録する(R11)。

【0076】 i 文字ハッシュテーブルに対する文書番号 $\#$ の登録処理が終了すると、インデックス i に1を加算して(R12)、R9へ戻り、加算された後のインデックス i が示すハッシュテーブルに対する文書番号 $\#$ の登録処理を開始する。

【0077】R9にて、加算後のインデックス*i*が使用ハッシュテーブル数*n*を越えると、今回読出した文字列に対する選択された全てのハッシュテーブル18～22に対する該当文書番号#の登録処理が終了したと判断して、R3へ戻り、先に取込んだ文書から次の文字列の読出を開始する。

【0078】また、検索処理部16が行う検索キーワード入力部15から入力された検索キーワードに対する検索処理を図6に示す流れ図を用いて説明する。検索キーワード入力部15から検索キーワードが入力されると(Q1)、入力検索キーワードの文字数*K*、及び該当文字列の文字種別を判断する。具体的には、前述した(a)、(b)、(c)に区分する(Q2)。そして、種別テーブル14から判別された文字種別に対応する使用ハッシュテーブルを特定する(Q3)。特定された使用ハッシュテーブルの数*n*と、使用ハッシュテーブルの名*Na*(1)、*Na*(2)、*Na*(3)、…、*Na*(*n*)と設定する(S4)。

【0079】以上の準備処理が終了すると、使用ハッシュテーブルを特定するインデックス*i*を1に初期化する(Q5)。そして、インデックス*i*が使用ハッシュテーブル数*n*以下で、かつ使用ハッシュテーブル名*Na*(*i*)が文字列の文字数*K*以下の場合(Q6)、*i*文字ハッシュテーブルに対する文書番号#の検索処理を開始する。

【0080】すなわち、該当文字列を*i*個の連続文字からなる複数の文字単位に分割して(Q7)、この各文字単位からハッシュ関数 $F_i(c_1, c_2, \dots, c_i)$ を用いて各ハッシュ値*h*を算出して、各ハッシュテーブル18～22のうちの*i*文字ハッシュテーブルの該当ハッシュ値*h*に対して設定されている各文書番号#を抽出(検索)する(Q8)。

【0081】*i*文字ハッシュテーブルに対する各文書番号#の検索処理が終了すると、インデックス*i*に1を加算して(Q9)、Q6へ戻り、加算された後のインデックス*i*に対応するハッシュテーブルに対する各文書番号#の検索処理を開始する。

【0082】Q6にて、加算後のインデックス*i*が使用ハッシュテーブル数*n*を越えると、今回入力した検索キーワードに対する選択された全てのハッシュテーブル18～22に対する文書番号#の検索処理が終了したと判断して、Q10へ進み、*n*個の各ハッシュテーブル18～22で検出された検索結果である全ての文書番号#のAND値を得る。具体的には、全ての文字単位が含まれる文書番号#を抽出する。そして、この軽策結果を検索結果出力部17へ表示出力する。

【0083】このように構成された第1実施形態の全文検索システムにおいては、登録する文字列や検索する検索キーワードの文字種別に応じて、登録したり検索に用いるハッシュテーブルの種別(*Na*(1)～*Na*(*n*))と数*n*とが異なる。

【0084】具体的には、文字種別(a)の英字、数字のみの場合は、1文字ハッシュテーブル18から5文字ハッシュテーブル22までの全てのハッシュテーブル5を使用する。逆に、文字種別(b)の英字、数字、ひらがな、カタカナ、漢字第一水準を含む場合は、1文字ハッシュテーブル18から3文字ハッシュテーブル20までの3つのハッシュテーブルを使用する。

【0085】したがって、4文字ハッシュテーブル21と5文字ハッシュテーブル22には文字種別(a)の英字、数字のみの組合せの文字単位に対するハッシュ値*h*のみしか設定されていない。英字、数字のみの組合せ数は文字種別(b)における組合せ数に比較して格段に少ないので、たとえこの4文字ハッシュテーブル21と5文字ハッシュテーブル22とをデータベース13に組込んだとしてもデータベース13の必要記憶容量が大幅に増加することはない。

【0086】次に、このように構成された第1実施形態の全文検索システムにおける具体的な登録動作及び検索動作を具体例を用いて説明する。先ず、次の2つの文書を登録する場合を説明する。

【0087】文書1(#=1) 「・・・12000件のデータ・・・」

文書2(#=2) 「・・・2000件中の120件・・・」

まず文書1(#=1)を読む。この場合、「12000」と「件のデータ」とに分割する。そして、同じ文字種別の文字を読むと、「12000」の部分が(a)の英数字なので、種別テーブル14で対応するハッシュテーブルを調べる。数字の場合は1～5文字ハッシュテーブル18～22が指定されているため、「12000」の文字列を次のように1～5文字の文字単位に分割して、それぞれのハッシュ値*h*を計算し、各ハッシュテーブル18～22に該当文書番号#(=1)を追加登録する。

【0088】

1文字ハッシュテーブル	1, 2, 0, 0, 0
2文字ハッシュテーブル	12, 20, 00, 00
3文字ハッシュテーブル	120, 200, 000
4文字ハッシュテーブル	1200, 2000
5文字ハッシュテーブル	12000

次の「件のデータ」の文字列は、ひらがな、カタカナ、漢字(第一水準)を含む(b)の文字種別である。この(b)の文字種別に対しては、種別テーブル14に1～3文字のハッシュテーブル18, 19, 20が指定されている。よって、同様な手法で各ハッシュテーブル18, 19, 20に対して該当文書番号#(=1)を追加登録する。

【0089】ここで、別の文字種別の文字列がつながる部分は少ない方の指定ハッシュテーブルに対して登録すると設定しておく、「00件のデータ」の中の次の文字に対して各ハッシュ値*h*を計算し、各ハッシュテーブ

ルに該当文書番号# (=1)を追加登録する。

【0090】

1文字ハッシュテーブル 件、の、デ、ー、タ

2文字ハッシュテーブル 0件、件の、のデ、デー、ー
タ

3文字ハッシュテーブル 00件、0件の、件のデ、の
デー、データ

次に、文書2 (#=2)を読む。同様に同じ文字種別の
文字を読み、次の通りに各ハッシュテーブルに該当文書
番号# (=2)を追加登録する。

【0091】1文字ハッシュテーブル 2, 0, 0,

0, 件、中、の、1, 2, 0

2文字ハッシュテーブル 20, 00, 00, 0件、件
中、中の、の1, 12, 20

3文字ハッシュテーブル 200, 000, 00件、0
件中、件中の、中の1, の12, 120

4文字ハッシュテーブル 2000

5文字ハッシュテーブル なし

以上で上記各文書 (#=1, #=2)の各ハッシュテー
ブルに対する登録が終了する。

【0092】次に、上述したように各文書 (#=1, #
=2)が登録されたデータベース13を検索キーワード
「12000」を用いて、この文字列「12000」を
含む文書を検索する場合を説明する。

【0093】まず検索キーワード「12000」が入力
されると、この検索キーワードの数字である文字種別
(a)を判定して、種別テーブル14から、文字種別(a)
に対応する1文字から5文字の各ハッシュテーブル18
~22を特定する。まず、検索キーワード「1200
0」を1文字単位に分割する。

【0094】1, 2, 0, 0, 0

これら5つの各1文字単位について、1文字ハッシュ
テーブル18に対して検索を行なうと、文書1、2が検索
される。次に検索キーワード「12000」を2文字単
位に分割する。

【0095】12, 20, 00, 00

これら4つの各2文字単位について、2文字ハッシュ
テーブル19に対して検索を行なうと、文書1、2が検索
される。次に検索キーワード「12000」を3文字単
位に分割する。

【0096】120, 200, 000

これら3つの各3文字単位について、3文字ハッシュ
テーブル20に対して検索を行なうと、文書1、2が検索
される。次に検索キーワード「12000」を4文字単
位に分割する。

【0097】1200, 2000

これら2つの各4文字単位について、4文字ハッシュ
テーブル21に対して検索を行なうと、文書1のみが検索
される。次に検索キーワード「12000」を5文字単
位に分割する。

【0098】12000

この1つの5文字単位について、5文字ハッシュテー
ブル22に対して検索を行なうと、文書1のみが検索され
る。

【0099】最後に、これまでの各検索結果のANDを
取ると文書1のみが残り、この文書1 (#=1)が最終
検索結果として出力される。これにより、正しい文書1
(#=1)のみが検索され、検索ノイズである文書2
(#=2)は検索されないため、検索精度が向上する。

【0100】このように、検索ノイズの発生確率が高い
英数字等の文字種別で構成される文字列に対するハッ
シュテーブルに対してのみ、構成文字数の高い4文字や5
文字のハッシュテーブル21, 22を設ければ、検索ノ
イズの発生確率が低下する。

【0101】なお、英数字の文字種は漢字に比較して格
段に少ないので、たとえ4文字や5文字のハッシュテー
ブル21, 22を設けたとしても、データベース13全
体の必要記憶容量が大幅に増加することはない。

【0102】(第2実施形態)図7は本発明の第2実施
形態に係わる全文検索システムの概略構成を示すブロッ
ク図である。図1に示す第1実施形態と同一部分には同
一符号が付してある。したがって、重複する部分の詳細
説明は省略されている。

【0103】この第2実施形態の全文検索システムのデ
ータベース13a内には、図1に示す第1実施形態と同
一の1文字ハッシュテーブル18、2文字ハッシュテー
ブル19、3文字ハッシュテーブル20の他に、共通ハ
ッシュテーブル23が設けられている。

【0104】この共通ハッシュテーブル23内には、図
9に示すように、構成する文字数が4文字、5文字、6
文字、…、n文字とそれぞれ異なる値を有する各文字
単位に対応する各文字組み合わせ毎に算出されたハッシュ
値hが登録されている。

【0105】具体的には、図示するように構成文字数i
毎に、ハッシュ関数 F_i が設定されている。例えば4文
字の場合はハッシュ関数 F_4 (c_1, c_2, c_3, c_4)
を用いてハッシュ値hを算出する。また、5文字の
場合はハッシュ関数 F_5 (c_1, c_2, c_3, c_4, c_5)
を用いてハッシュ値hを算出する。さらに、6文字
の場合はハッシュ関数 F_6 ($c_1, c_2, c_3, c_4, c_5, c_6$)
を用いてハッシュ値hを算出する。

【0106】また、種別テーブル14a内には、検索キ
ーワード及び登録文字列の前述した(a), (b), (c)の
文字種別毎に利用する各ハッシュテーブルが登録されて
いる。

【0107】(a)に示す英字、数字のみの場合、1文字
ハッシュテーブル18、2文字ハッシュテーブル19、
3文字ハッシュテーブル20及び共通ハッシュテーブル
23が登録されている。なお、文字種別(b), (c)に
対しては図2で示した第1実施形態の種別テーブル14と

同一のハッシュテーブルが設定されている。

【0108】このような構成の第2実施形態の全文検索システムにおいて、登録処理部12は登録文書入力部11から入力された各文書のそれぞれ登録すべき各文字列を例えば1～N個の文字からなる各単位文字に分離して、種別テーブル14aの指定する各ハッシュテーブルへ該当文章番号#を登録するが、4文字以上の文字単位に対する文書番号#の登録は全て共通ハッシュテーブル23へ一括して登録される。

【0109】検索処理部16は、検索キーワード検索部15から入力された検索キーワードでデータベース13aを検索する場合においても、検索キーワードを1～N個の文字からなる各文字単位に分離して、種別テーブル14aの指定する各ハッシュテーブルを検索するが、4文字以上の文字単位に対する検索はすべて共通ハッシュテーブル23に対して実施する。

【0110】このような、共通ハッシュテーブル23を使用したとしても、目標とする文書を確実に検索できる。また、図9に示したように、ハッシュ値hは各構成文字数毎のハッシュ関数 F_4 、 F_5 、 F_6 、…、 F_n 毎に個別の値として求まるが、たとえ構成文字数が異なる文字列であっても、4文字単位、5文字単位、6文字単位等の構成文字数が大きいものは登録されている各文書における発生確率が小さい。したがって、この各文字数毎のハッシュ値hを一つの共通ハッシュテーブル23に組込んだとしても誤った文書が検索される検索ノイズの発生確率は大幅に上昇することはない。

【0111】よって、検索精度を低下させることなく、データベース13aの必要記憶容量を節減できる。

(第3実施形態) 図10は本発明の第3実施形態に係わる全文検索システムの概略構成を示すブロック図である。図1に示す第1実施形態と同一部分には同一符号が付してある。したがって、重複する部分の詳細説明は省略されている。

【0112】この第3実施形態の全文検索システムのデータベース13b内には、1文字ハッシュテーブル18a、2文字ハッシュテーブル19a、3文字ハッシュテーブル20a、4文字ハッシュテーブル21a、5文字ハッシュテーブル22a、6文字ハッシュテーブル24、7文字ハッシュテーブル25の合計7個のハッシュテーブルが設けられている。

【0113】また、種別テーブル14b内には、図11に示すように、検索キーワード及び登録文字列の(a)、(b1)、(b2)、(c)の合計4種類の文字種別毎に利用する各ハッシュテーブル及び採用する各ハッシュ関数が登録されている。

【0114】(a)の文字種別は、第1実施形態と同様に英字と数字のみであり、この文字種別(a)に対して前述した1文字から7文字までの全てのハッシュテーブル18a～25が使用ハッシュテーブルとして登録されて

いる。さらに、各ハッシュテーブル18a～25毎に採用するハッシュ関数 F_1 、 F_2 、 F_3 、 F_4 、 F_5 、 F_6 、 F_7 が登録されている。各ハッシュ関数 F_i ～ F_7 は文字単位の構成文字数iが異なるのみのである同一種類のハッシュ関数 F_i (c_1 、…、 c_i)である。

【0115】(b1)の文字種別は、ひらかなとカタカナのみであり、この文字種別(b1)に対して前述した1文字から3文字までの各ハッシュテーブル18a～20aが使用ハッシュテーブルとして登録されている。さらに、各ハッシュテーブル18a～20a毎に採用するハッシュ関数 G_1 、 G_2 、 G_3 が登録されている。各ハッシュ関数 G_1 ～ G_3 は文字単位の構成文字数iが異なるのみの度同一種類のハッシュ関数 G_i (c_1 、…、 c_i)である。

【0116】(b2)の文字種別は、第一水準の漢字のみであり、この文字種別(b2)に対して前述した1文字から3文字までの各ハッシュテーブル18a～20aが使用ハッシュテーブルとして登録されている。さらに、各ハッシュテーブル18a～20a毎に採用するハッシュ関数 D_1 、 D_2 、 D_3 が登録されている。各ハッシュ関数 D_1 ～ D_3 は文字単位の構成文字数iが異なるのみの同一種類のハッシュ関数 D_i (c_1 、…、 c_i)である。

【0117】(c)の文字種別は、第1実施形態と同様に第一水準の漢字と外字のみであり、この文字種別(c)に対して前述した1文字ハッシュテーブル18aと2文字ハッシュテーブル19aが使用ハッシュテーブルとして登録されている。さらに、各ハッシュテーブル18a、19aにそれぞれ採用するハッシュ関数 E_1 、 E_2 が登録されている。各ハッシュ関数 E_1 、 E_2 は文字単位の構成文字数iが異なるのみの同一種類のハッシュ関数 E_i (c_1 、…、 c_i)である。

【0118】このように、検索キーワード及び登録文字列の前述した(a)、(b1)、(b2)、(c)の合計4種類の文字種別毎に異なる種類のハッシュ関数 F_i 、 G_i 、 D_i 、 E_i が設定されている。

【0119】登録処理部12aは、登録文書入力部11から入力された各文書の各登録すべき各文字列の文字種別(a)、(b1)、(b2)、(c)を判断して、例えば1～N個の文字からなる各単位文字に分離する。そして、種別テーブル14bの該当文字種別に指定されたハッシュ関数 F_i 、 G_i 、 D_i 、 E_i を用いてハッシュ値hを算出し、同じく種別テーブル14bで指定されたハッシュテーブルの該当ハッシュ値hの欄に今回登録しようとする文字列が含まれる文書番号#を追加登録する。

【0120】検索処理部16aは、検索キーワード検索部15から入力された検索キーワードでデータベース13bを検索する場合、検索キーワードを構成する文字の文字種別(a)、(b1)、(b2)、(c)を判断して、例えば1～N個の文字からなる各単位文字に分割する。そして、種別テーブル14bの該当文字種別に指定されたハッシ

関数 F_i , G_i , D_i , E_i を用いてハッシュ値 h を算出し、同じく種別テーブル 14 b で指定されたハッシュテーブルの該当ハッシュ値 h の欄に登録された文書番号を読み取る。

【0121】このような、検索キーワード、登録する文字列を構成する各文字の文字種別(a) , (b1) , (b2) , (c) 毎に異なるハッシュ関数 F_i , G_i , D_i , E_i を用いてハッシュ値 h を算出したとしても、目標とする文書を確実に検索できる。

【0122】さらに、この第3実施形態においては、各文字種別(a) , (b1) , (b2) , (c) 毎に異なるハッシュ関数 F_i , G_i , D_i , E_i を用いてハッシュ値 h を算出している。

【0123】一般に、各文字種別毎に、1文字単位、2文字単位、3文字単位の各文字の組合せの発生状況が異なるので、全ての文字種別(a) , (b1) , (b2) , (c) に亘って同一種別のハッシュ関数を採用してハッシュ値 h を算出した場合においては、文字種別によっては、ハッシュテーブル内において、ある特定のハッシュ値 h に対して多くの文書番号が登録されることになる。その結果、検索ノイズの発生確率が上昇したり、ハッシュテーブルを有効に使用できない懸念がある。

【0124】したがって、各文字種別(a) , (b1) , (b2) , (c) 毎に、該当文字種別の組合せの発生状況に対応した最適のハッシュ関数 F_i , G_i , D_i , E_i を設定することによって、一つのハッシュ値 h に対して多数の文書番号が登録されることを抑制でき、検索ノイズの発生確率を低下でき、検索精度を向上できる。

【0125】また、4文字から7文字までの各ハッシュテーブル 22 a , 24 , 25 内には、英字及び数字からなる文字単位の組合せに対するハッシュ値 h しか登録されていない。この英字及び数字からなる文字単位の組合せ数は、漢字の組合せ数の比較して格段に小さいので、同一の組合せが同一ハッシュ値 h にならないように、この文字種別のハッシュ関数 F_4 , F_5 , F_6 , F_7 を調整することによって、たとえ検索キーワードに数字が多く含まれる場合であっても、検索ノイズの発生確率を低下でき、検索精度を向上できる。

【0126】なお、本発明は上述した各実施形態のみに限定されるものではない。例えば図1に示した全文検索システムの登録処理部12、検索処理部16の機能をプログラム化し、予めCD-ROMなどの記録媒体に書き込んでおき、このCD-ROMをCD-ROMドライブを搭載した計算機に装填し、計算機がCD-ROMからプログラムをロードすることにより上記実施形態と同様の機能を実現することができる。なお、記録媒体としては、上記CD-ROM以外に、磁気テープ、DVD-ROM、フロッピーディスク、MO MD、CD-R、メモ리카ードなどでもよい。

【0127】

【発明の効果】以上説明したように本発明の全文検索システム及び全文検索方法並びに全文検索プログラムを記録した記録媒体においては、検索情報テーブルの設定数を検索すべき文字列の文字種別に応じた値に設定している。したがって、検索情報テーブルを記憶するデータベースの記憶容量を大幅に増加することなく、検索精度を向上できる。

【0128】また、文字数が互いに異なる複数種類の文字列に対応するハッシュ値が組込まれた検索情報テーブルを用いるので、検索精度を低下することなく検索情報テーブルを記憶するデータベースの記憶容量を低減できる。

【0129】さらに、文字列を構成する文字の種別毎に異なるハッシュ関数を用いて算出されたハッシュ値を検索情報テーブルに設定している。したがって、たとえば、発生確率の高い文字単位どうしを同一ハッシュ値に設定されることを防止するように、各文字種別毎にハッシュ関数を設定することによって、検索ノイズの発生確率を低減できる。

【図面の簡単な説明】

【図1】 本発明の第1実施形態に係わる全文検索方法を採用した全文検索システムの概略構成を示すブロック図

【図2】 同全文検索システムに組込まれた種別テーブルの登録内容を示す図

【図3】 同全文検索システムのデータベースに組込まれた各ハッシュテーブルの登録内容を示す図

【図4】 同じく同データベースに組込まれた各ハッシュテーブルの登録内容を示す図

【図5】 同全文検索システムのデータベースに対する文書の登録処理を示す流れ図

【図6】 同全文検索システムのデータベースに対する文書の検索処理を示す流れ図

【図7】 本発明の第2実施形態に係わる全文検索システムの概略構成を示すブロック図

【図8】 同全文検索システムに組込まれた種別テーブルの登録内容を示す図

【図9】 同全文検索システムのデータベースに組込まれた共通ハッシュテーブルの登録内容を示す図

【図10】 本発明の第3実施形態に係わる全文検索システムの概略構成を示すブロック図

【図11】 同全文検索システムに組込まれた種別テーブルの登録内容を示す図

【図12】 従来の全文検索システムの概略構成を示すブロック図

【図13】 同全文検索システムに組込まれた各文字検索情報テーブルの記憶内容を示す図

【図14】 同全文検索システムデータベースに対する文書の検索処理を示す流れ図

【図15】 検索キーワードを文字単位に分割する場合

の分割種別を示す図

【図16】 一般的なi文字ハッシュテーブルの記憶内容を示す図

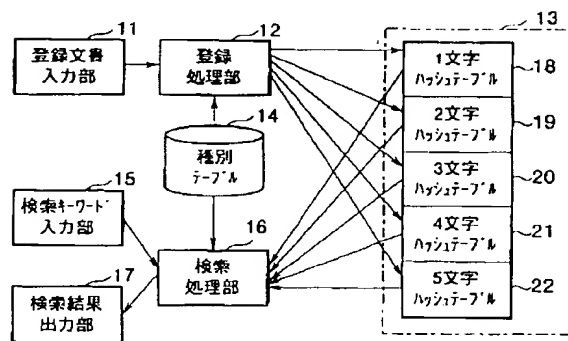
【符号の説明】

- 11…登録文書入力部
12、12a…登録処理部
13、13a、13b…データベース
14、14a、14b…種別テーブル
15…検索キーワード入力部
16、16a…検索処理部
17…検索結果出力部

17…検索結果出力部

- 18、18a…1文字ハッシュテーブル
19、19a…2文字ハッシュテーブル
20、20a…3文字ハッシュテーブル
21、21a…4文字ハッシュテーブル
22、22a…5文字ハッシュテーブル
23…共通ハッシュテーブル
24…6文字ハッシュテーブル
25…7文字ハッシュテーブル

【図1】



【図2】

種別テーブル		
キーコード、文字列	種別	対応ハッシュテーブル
a	英字、数字のみ	1, 2, 3, 4, 5文字ハッシュテーブル
b	ひらがな、カタカナ漢字第一水準、英字、数字を含む	1, 2, 3文字ハッシュテーブル
c	漢字第二水準、外字のみ	1, 2文字ハッシュテーブル

【図3】

【図4】

(a)

4文字ハッシュテーブル										
文書# ハッシュ値h	1	2	3	4	5	n
1	0	0	0	0	1					1
2	0	1	0	0	0					0
3	1	0	1	0	0					1
...										
N4	0	0	0	0	1					0

(b)

5文字ハッシュテーブル										
文書# ハッシュ値h	1	2	3	4	5	n
1	0	0	1	0	0					0
2	0	0	0	0	1					0
...										
N5	0	0	0	1	0					1

(a)

1文字ハッシュテーブル										
文書# ハッシュ値h	1	2	3	4	5	n
1	0	0	1	0	1					0
2	0	1	0	1	0					1
3	0	0	0	1	1					0
...										
N1	0	1	0	0	0					0

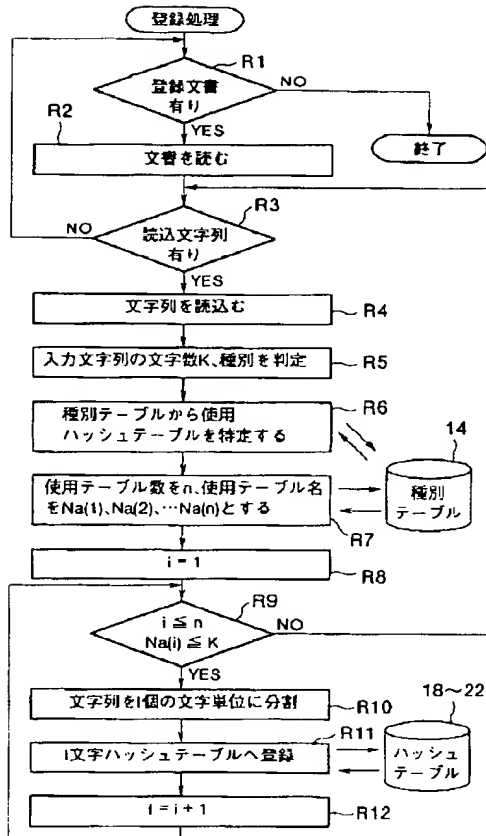
(b)

2文字ハッシュテーブル										
文書# ハッシュ値h	1	2	3	4	5	n
1	0	1	0	0	1					1
2	0	1	0	0	1					0
3	1	0	1	0	0					1
...										
N2	0	0	1	0	1					0

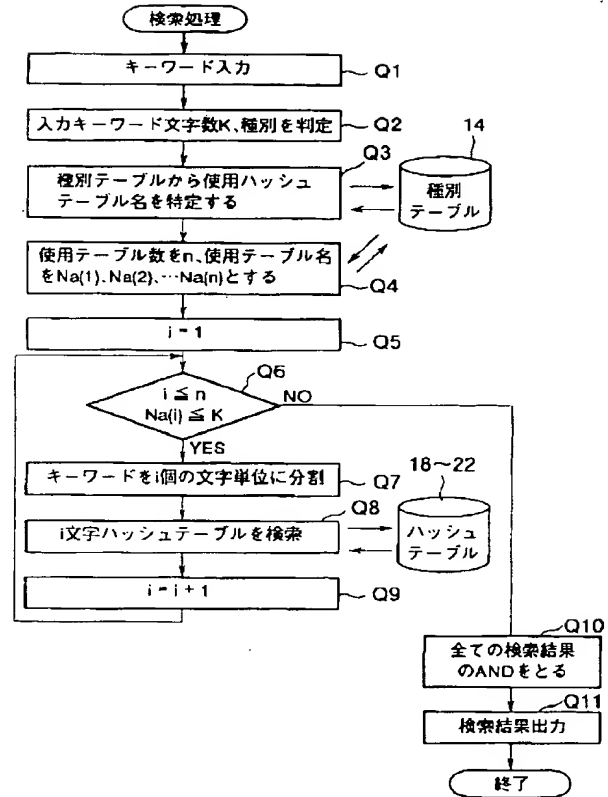
(c)

3文字ハッシュテーブル										
文書# ハッシュ値h	1	2	3	4	5	n
1	0	0	1	0	0					0
2	0	1	0	0	1					0
...										
N3	0	1	0	1	0					1

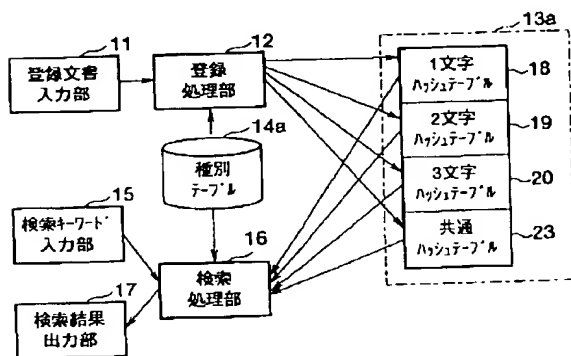
【図5】



【図6】



【図7】



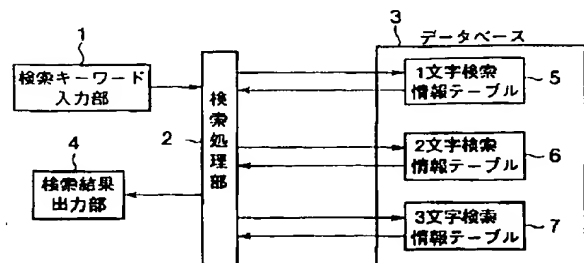
【図8】

	キーワード、文字列	種別	対応ハッシュテーブル
a	英字、数字のみ		1、2、3文字、共通ハッシュテーブル
b	ひらがな、カタカナ漢字第一水準を含む		1、2、3文字ハッシュテーブル
c	漢字第二水準、外字のみ		1、2文字ハッシュテーブル

【図12】

【図15】

- 1文字単位 日、本、の、技、術
 2文字単位 日本、本の、の技、技術、
 3文字単位 日本_の、本の技、の技術、
- 【日本の技術】



【図9】

文書 #	1	2	3	4	5	...	n
ハッシュ値 h	0	0	1	0	0		0
1	0	0	1	0	0		0
2	0	1	0	0	1		0
...							
Nk	0	1	0	1	0		1

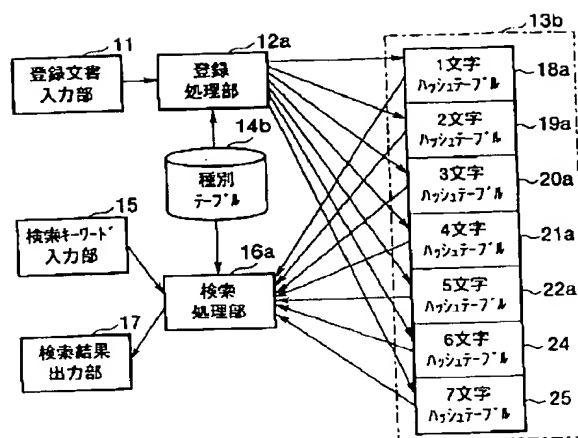
ハッシュ値 h

4文字の時 $F_4(C_1, C_2, \dots, C_4)$
 5文字の時 $F_5(C_1, C_2, \dots, C_5)$
 6文字の時 $F_6(C_1, C_2, \dots, C_6)$
 ...
 n文字の時 $F_n(C_1, C_2, \dots, C_n)$

F_i : i文字ハッシュ関数 ($i = 4, \dots, n$)

C_j : j番目の文字 ($j = 1, \dots, n$)

【図10】



【図11】

	ワード、文字列	種別	対応ハッシュテーブル	ハッシュ関数
a	英字 数字	1	1	$F_1(C_1)$
		2	2	$F_2(C_1, C_2)$
		3	3	$F_3(C_1, C_2, C_3)$
		4	4	$F_4(C_1, \dots, C_4)$
		5	5	$F_5(C_1, \dots, C_5)$
		6	6	$F_6(C_1, \dots, C_6)$
		7	7	$F_7(C_1, \dots, C_7)$
b1	ひらがな カタカナ	1	1	$G_1(C_1)$
		2	2	$G_2(C_1, C_2)$
		3	3	$G_3(C_1, C_2, C_3)$
b2	漢字第一水準	1	1	$D_1(C_1)$
		2	2	$D_2(C_1, C_2)$
		3	3	$D_3(C_1, C_2, C_3)$
c	漢字第二水準、外字	1	1	$E_1(C_1)$
		2	2	$E_2(C_1, C_2)$

F_i : ハッシュ関数

G_i : ハッシュ関数

D_i : ハッシュ関数

E_i : ハッシュ関数

C_j : j番目の文字列

【図13】

文書 #	1	2	3	4	...	n
文字単位	0	1	1	1	0	0
あ	0	1	1	1	0	0
い						
...						
東	1	0	1	0	0	1

(a)

文書 #	1	2	3	4	...	n
文字単位	0	0	1	0	0	0
ああ	0	0	1	0	0	0
あい						
...						
東京	1	0	0	0	0	1

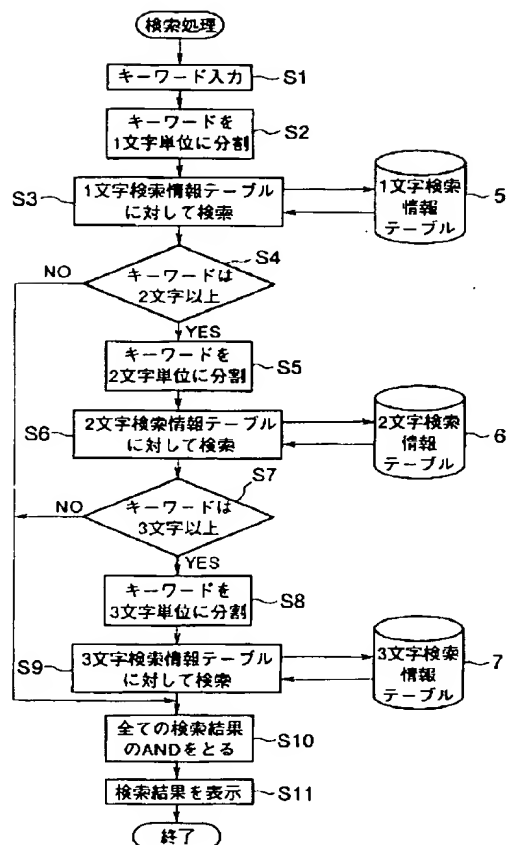
(b)

【図16】

文書 #	1	2	3	4	5	...	n
ハッシュ値 h	0	0	1	0	1		0
1	0	0	1	0	1		1
2	0	1	0	1	0		0
3	0	0	0	1	1		1
...							
Ni	0	1	0	0	0		0

h = $F(C_1, C_2, \dots, C_i)$

【図14】



フロントページの続き

(72)発明者 杉山 晋也
東京都府中市東芝町1番地 株式会社東芝
府中工場内

(72)発明者 菅谷 友秀
東京都府中市東芝町1番地 株式会社東芝
府中工場内